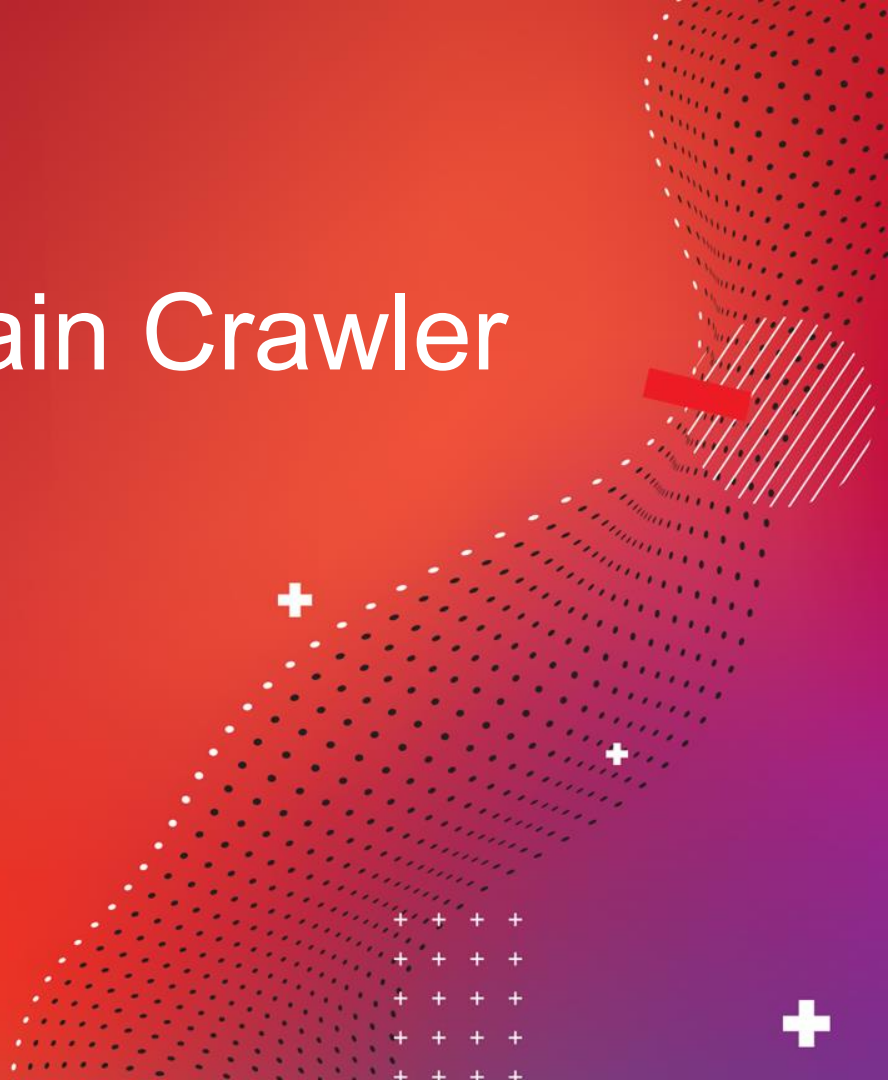


Пишем свой Domain Crawler

Евгений Карагодин



HighLoad++
Весна 2021



Евгений Карагодин



Занимаюсь разработкой 13 лет

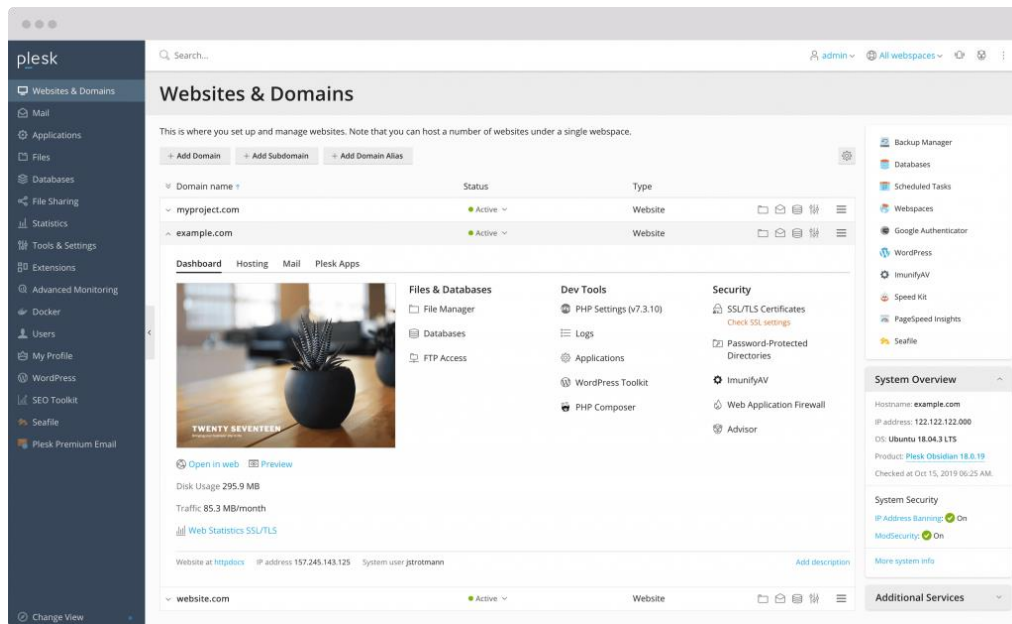


Люблю обучать и делиться
знаниями



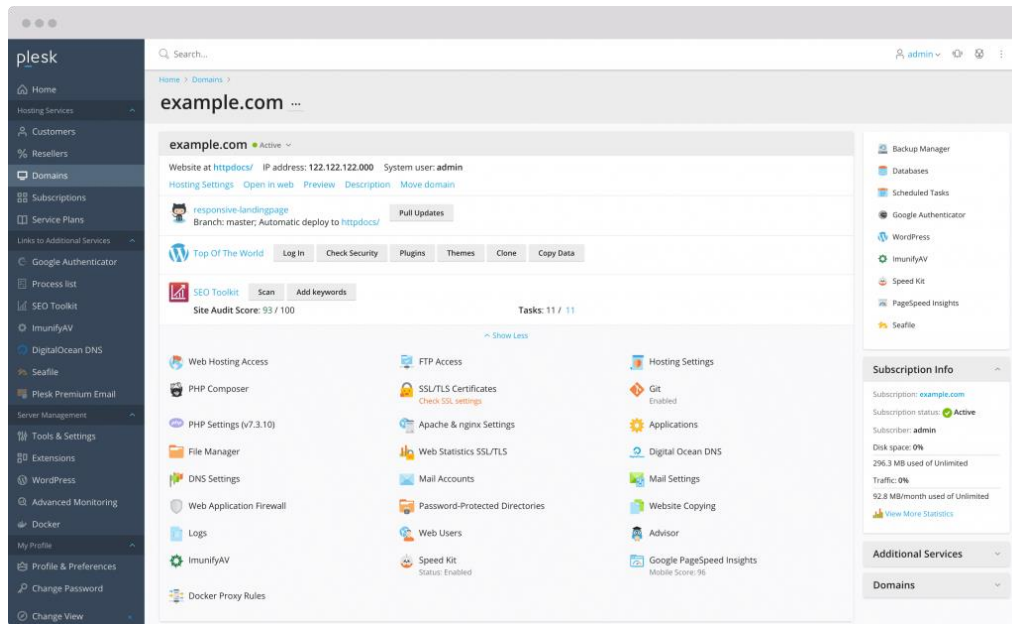
Plesk

Основной продукт — панель для управления сервером



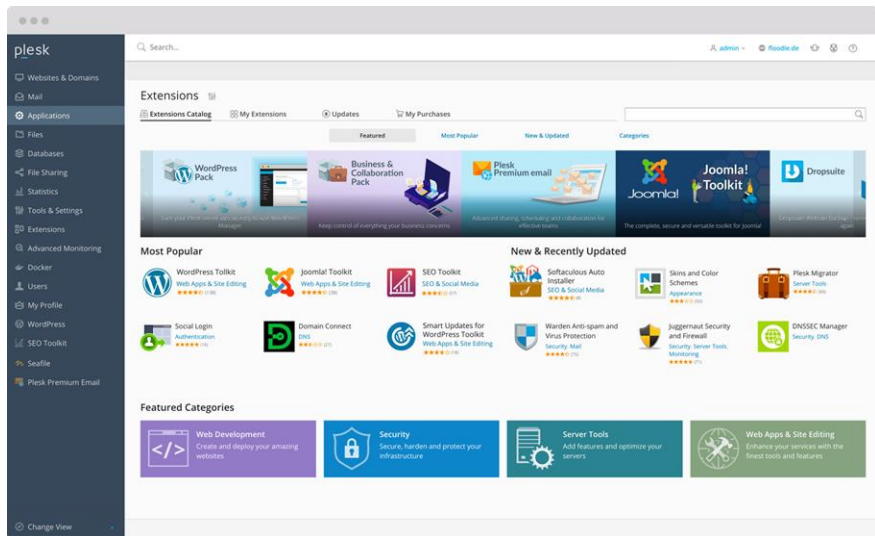
Plesk

Обслуживаем 11 миллионов сайтов — это ≈6% Интернет



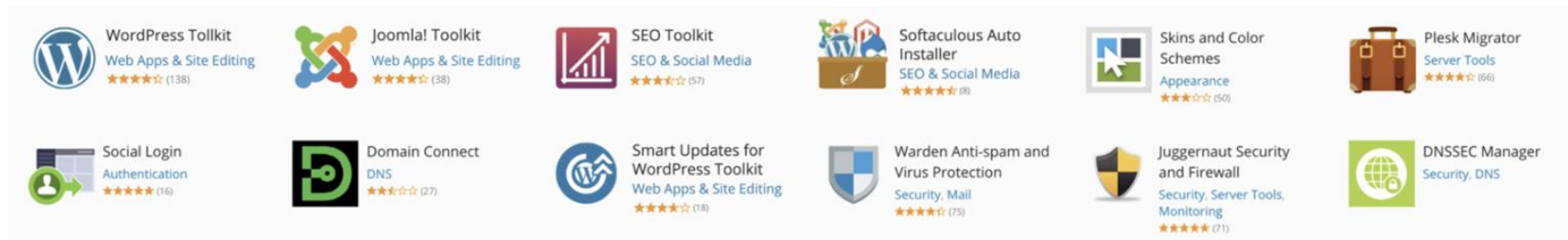
Plesk

- Мы хотим знать, какие технологии популярны у наших клиентов
- Что популярно в Web в целом



Что мы хотим знать?

- Использование CDN — нужно собирать IP-адреса
- Какие использует nameservers — нужно собирать DNS-записи
- Данные по SSL-сертификатам: валидность, тип, кем выпущен
- Веб-сервер, язык программирования, CMS, ...



Как это можно узнать?

- Аналитика из продукта
- Купить готовые отчёты
- Собирать данные с сайтов

Аналитика из продукта

- Не всё можно узнать изнутри
- Хотим знать не только про наших клиентов
- Делать отдельный сервис проще, чем основной продукт

Купить готовые отчёты

- Стоит \$18000 в год
- Методика сбора данных не раскрывается
- Сложно проверить качество данных
- Невозможно выделить данные наших клиентов



Придётся парсить сайты!



Как парсить сайты?

- Берём список доменов
- Делаем запрос на каждый домен
- Сохраняем результат

как парсить сайты



All



Videos



Images



News



Maps



More

Tools

About 197,000 results (0.50 seconds)

Это вообще законно?

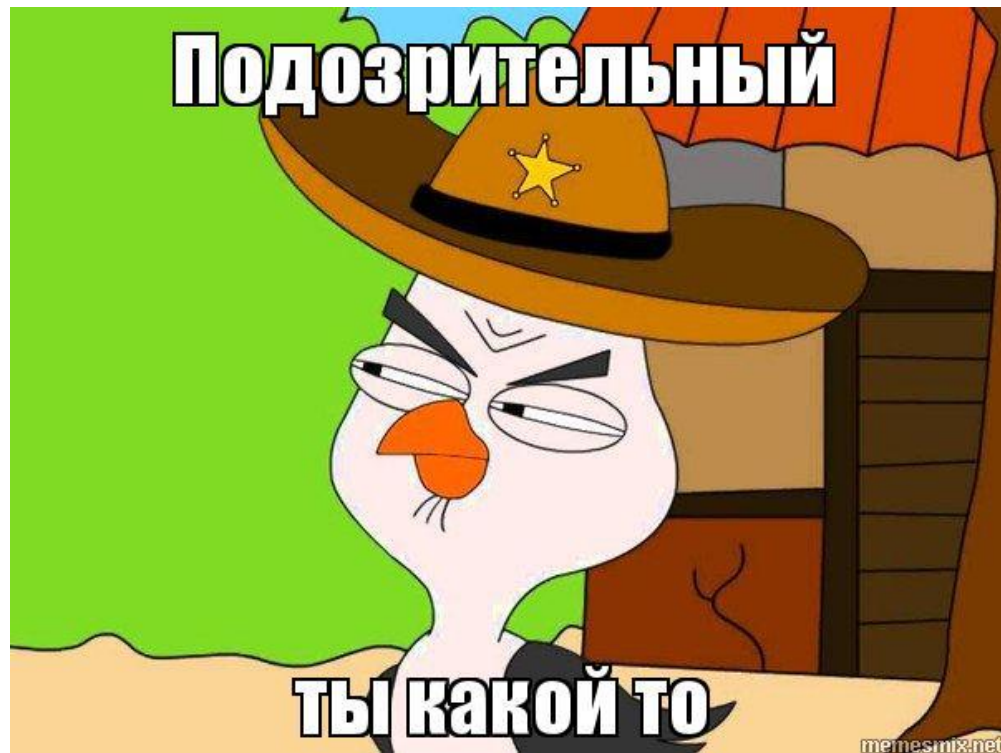


Я вызываю полицию!

Первое правило

Не маскируйся

User-agent: PleskBot



Второе правило

Уважай robots.txt

User-agent: *

Disallow: /

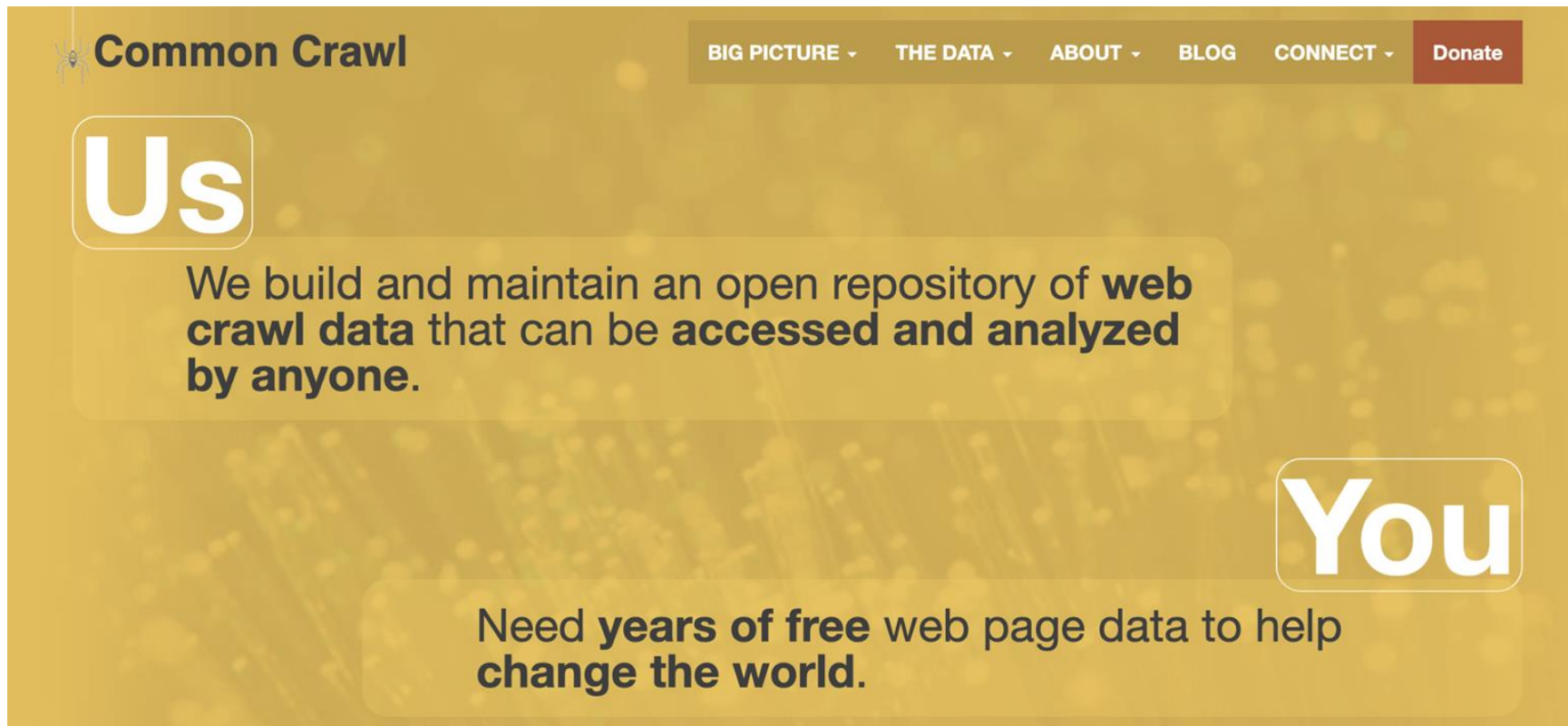


Третье правило

Не наноси урон сайту

- Три простых правила
- Учитывайте лицензию, по которой опубликован контент
- Учитывайте локальные законы
- Иначе вас будут банить, вам грозят серьёзные штрафы и репутационные потери

Как парсить сайты, не создавая нагрузку на них?



The image shows the header and main content of the Common Crawl website. The header is a dark blue bar with the Common Crawl logo (a spider) and the text 'Common Crawl' on the left. On the right, there are navigation links: 'BIG PICTURE', 'THE DATA', 'ABOUT', 'BLOG', 'CONNECT', and a red 'Donate' button. The main content area has a light blue background with a subtle pattern of yellow dots and lines. It features a large 'Us' in a rounded rectangle, followed by the text 'We build and maintain an open repository of **web crawl data** that can be **accessed and analyzed** by anyone.' To the right, there is a large 'You' in a rounded rectangle, followed by the text 'Need **years of free** web page data to help **change the world.**'

Common Crawl

[BIG PICTURE](#) [THE DATA](#) [ABOUT](#) [BLOG](#) [CONNECT](#) [Donate](#)

Us

We build and maintain an open repository of **web crawl data** that can be **accessed and analyzed** by anyone.

You

Need **years of free** web page data to help **change the world.**

Common Crawl

- Примерно 90 миллионов доменов
- Пересечение с нашим списком — 27 миллионов
- Нет нужных данных

Чем парсить?

- Наш выбор Node.js.
- Скриптовый язык — удобен для быстрого прототипирования.
- Асинхронный ввод/вывод — большая часть времени уходит на сетевые задержки.
- Большая экосистема npm-пакетов — есть готовые решения на все случаи.
- Если проблема возникает на уровне С-кода, то с этим ничего нельзя поделать.

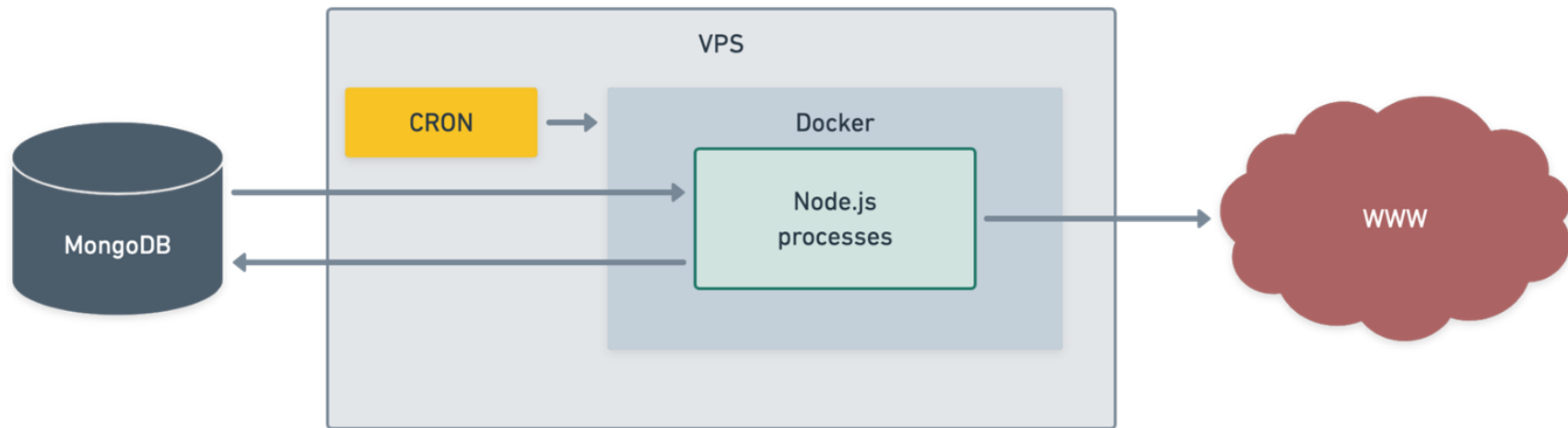
Node.js vs Golang

- На Go будет быстрее!!!1
- Нет, потому что дольше ждём сеть

MVP

- Несколько тысяч доменов в сутки
- Три скрипта на node.js
- Запускаются по крону друг за другом
- Данные хранятся либо в памяти, либо в локальных файлах
- Всё работает в докере

MVP

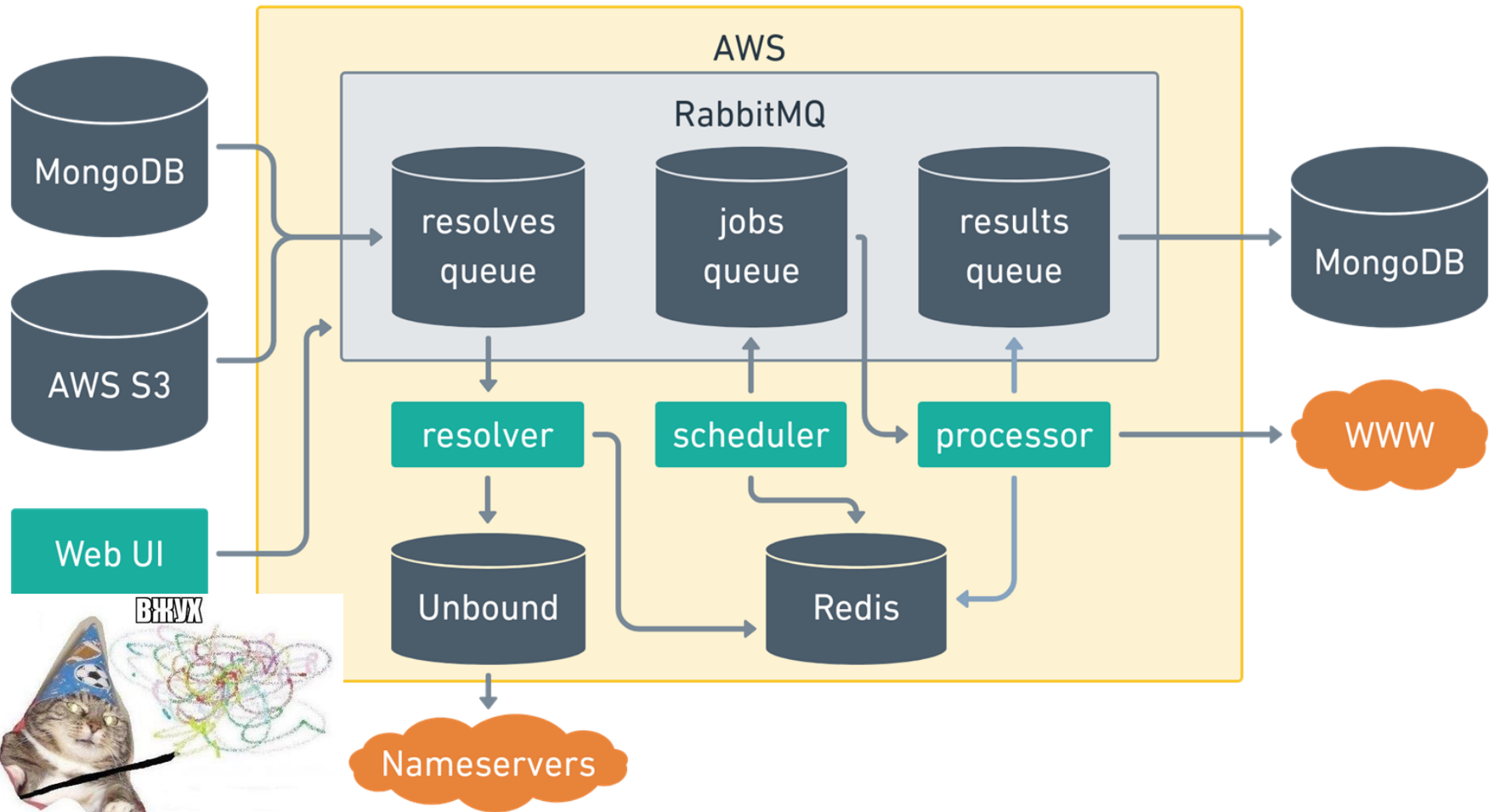


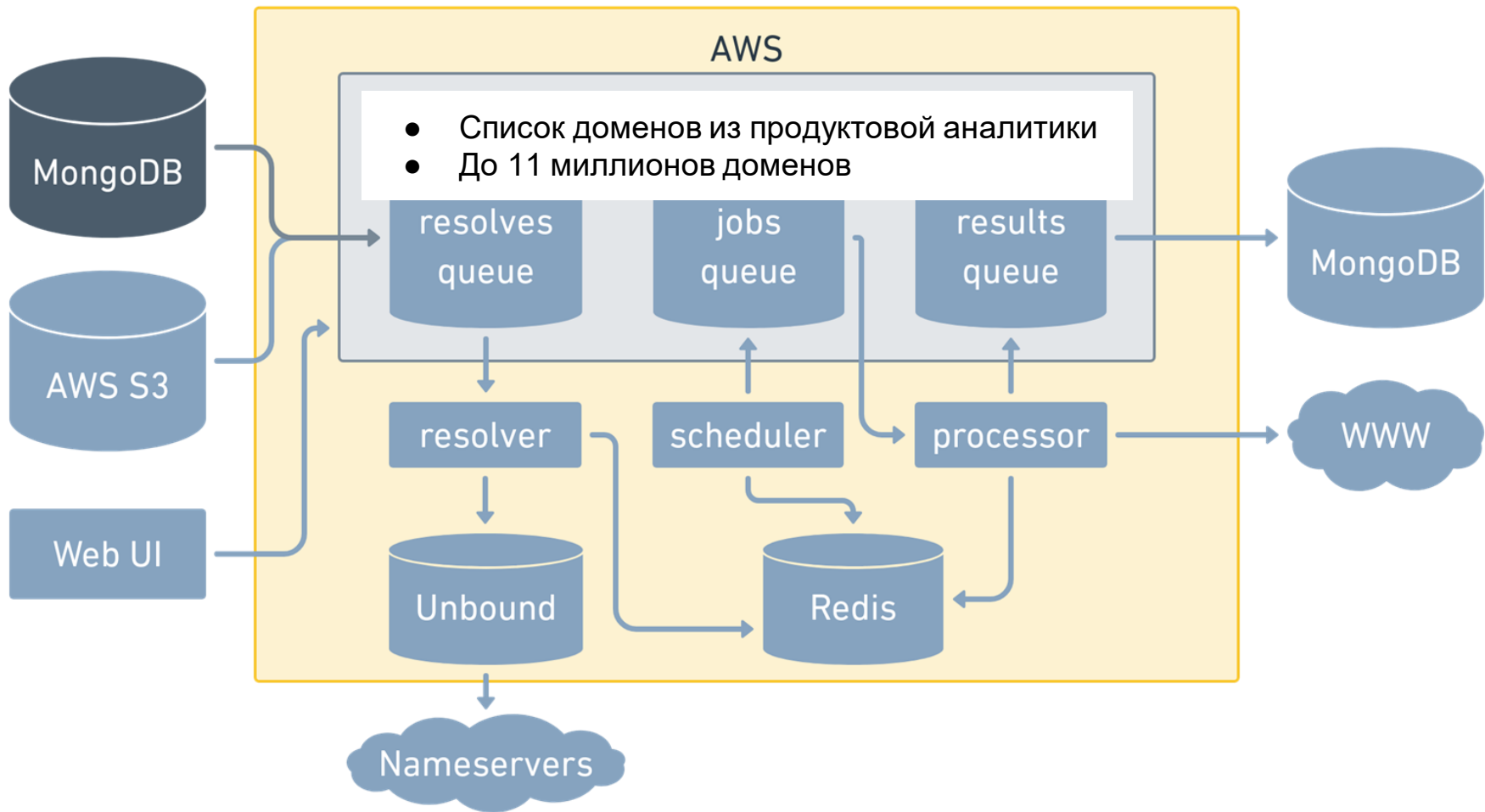
Результат

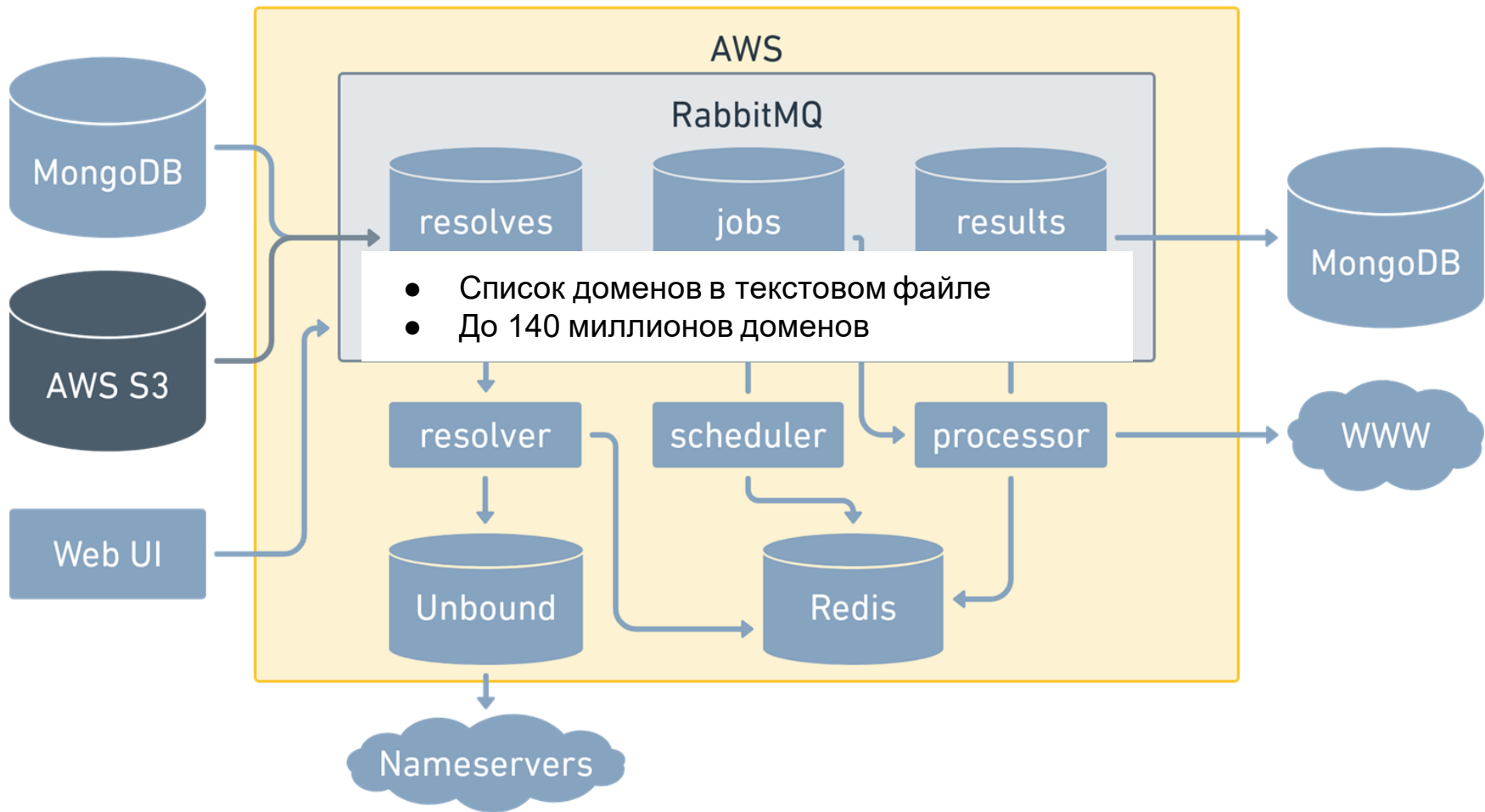
- Сделали быстро
- Получили результат
- Запланировали обходить десятки миллионов доменов на регулярной основе и сотни миллионов по требованию

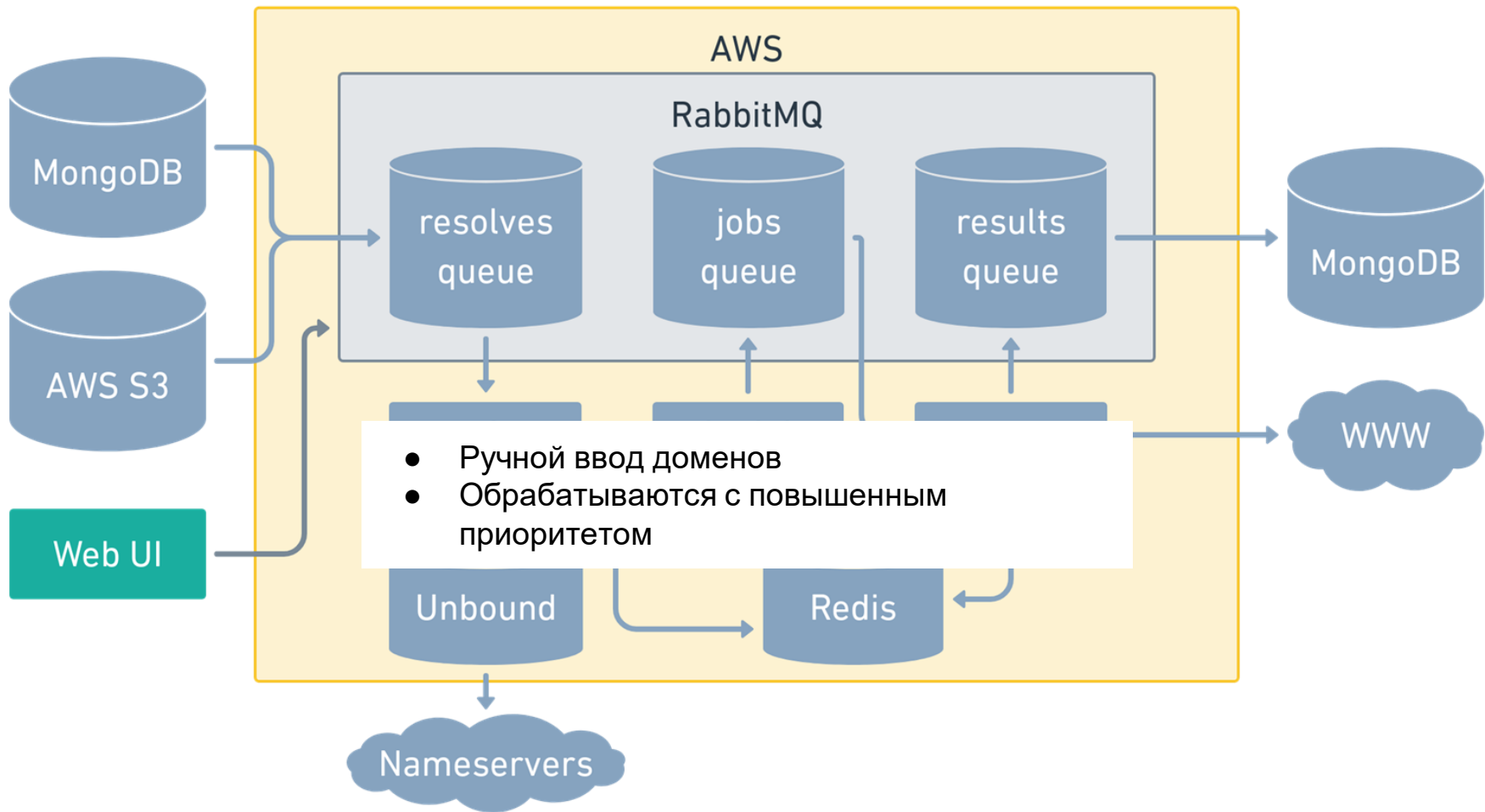
Требования Пожелания

- Масштабирование
- Надёжность
- Минимум ресурсов на разработку







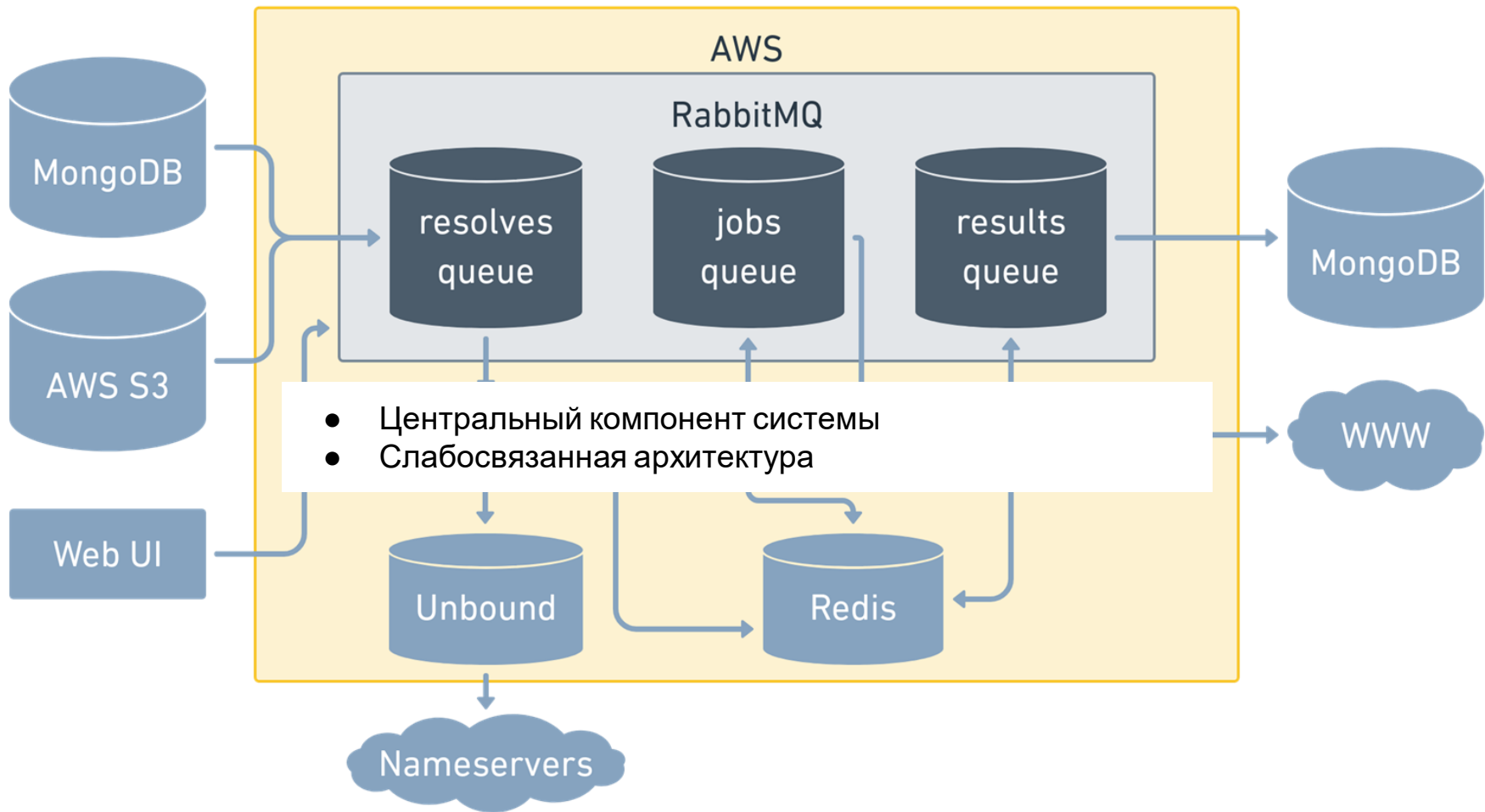


Domain name *

https://www.highload.ru/

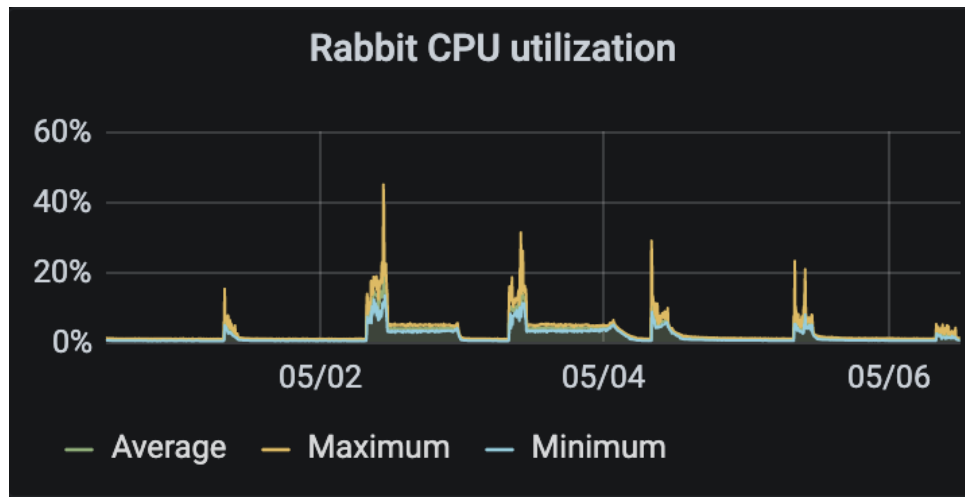
OK

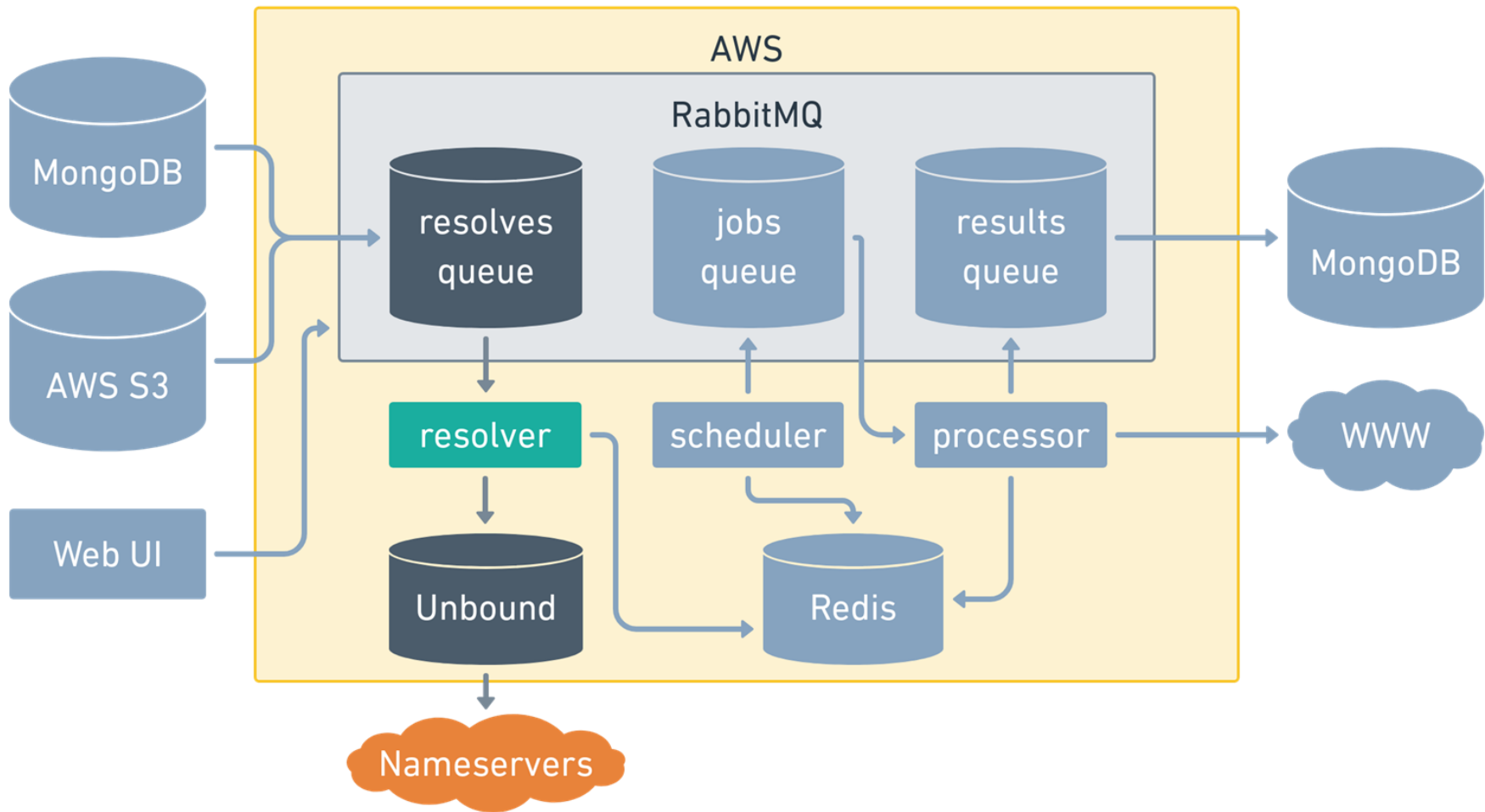
```
{
  "domain": "www.highload.ru",
  "metaInfo": {
    "domainType": "domain",
    "domainNameSource": "web ui"
  },
  "resolve": {
    "ip": "178.248.233.16"
  },
  "resolveNs": [
    {
      "domain": "ns2.reg.ru",
      "ip": "176.99.13.12"
    },
    {
      "domain": "ns1.reg.ru",
      "ip": "194.58.117.15"
    }
  ],
  "robotsTxt": {
    "allowed": true
  },
  "http": {
    "responseHeaders": {
      "server": "QRATOR",
      "date": "Thu, 06 May 2021 04:03:39 GMT",
      "content-type": "text/html",
```



Очередь

- Попробовали Bull (очередь на основе Redis)
- Bull упирался в CPU
- Поменяли Bull на RabbitMQ
- Большая производительность на том же железе



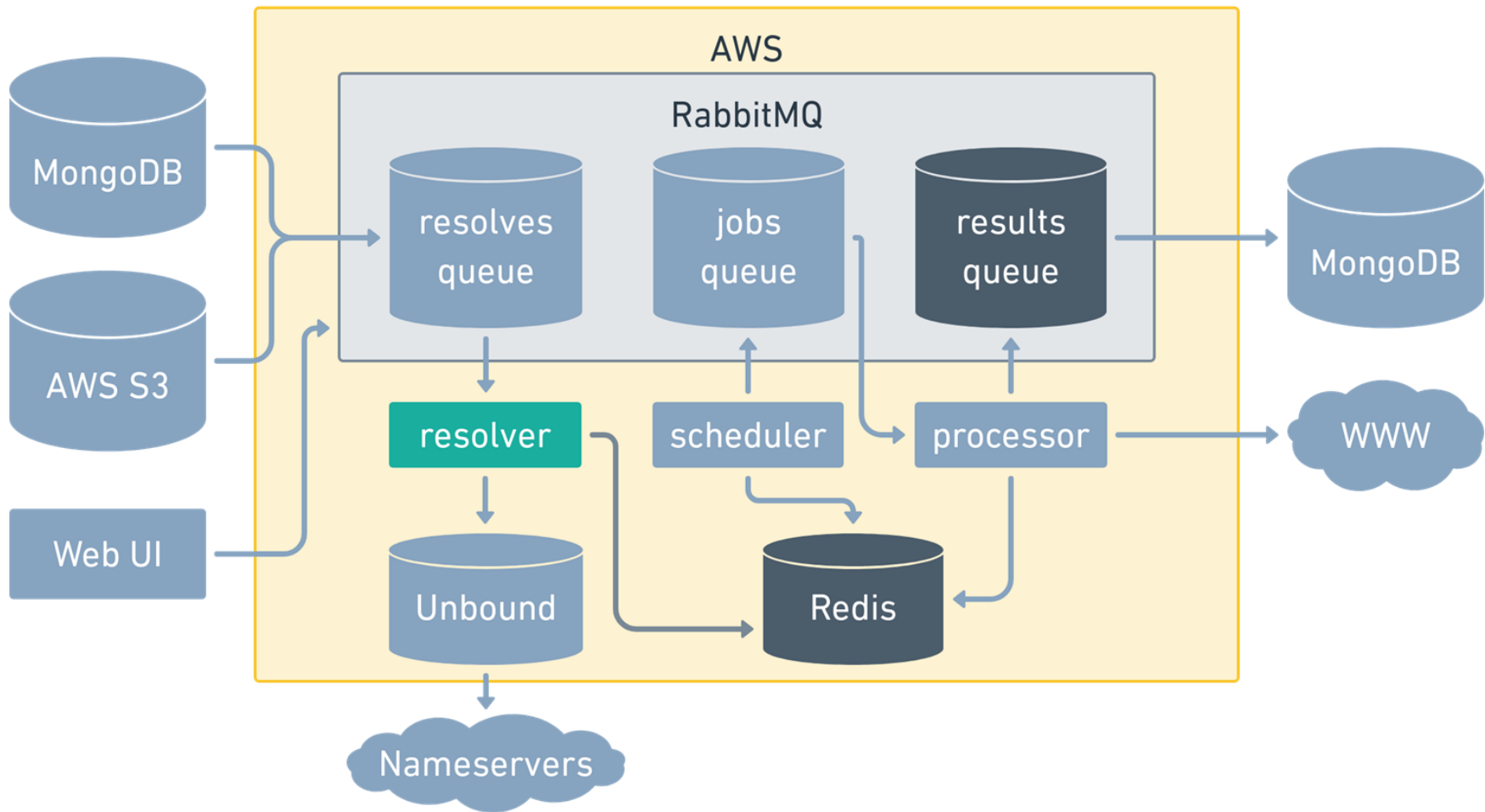


Массированный резолвинг доменов

- Не все резолверы выдержат
- Положили локальный DNS-резолвер
- Cloudflare резолвит домены, но не отдаёт все записи
- Подняли Unbound для кэширования
- ≈ 1000 rps

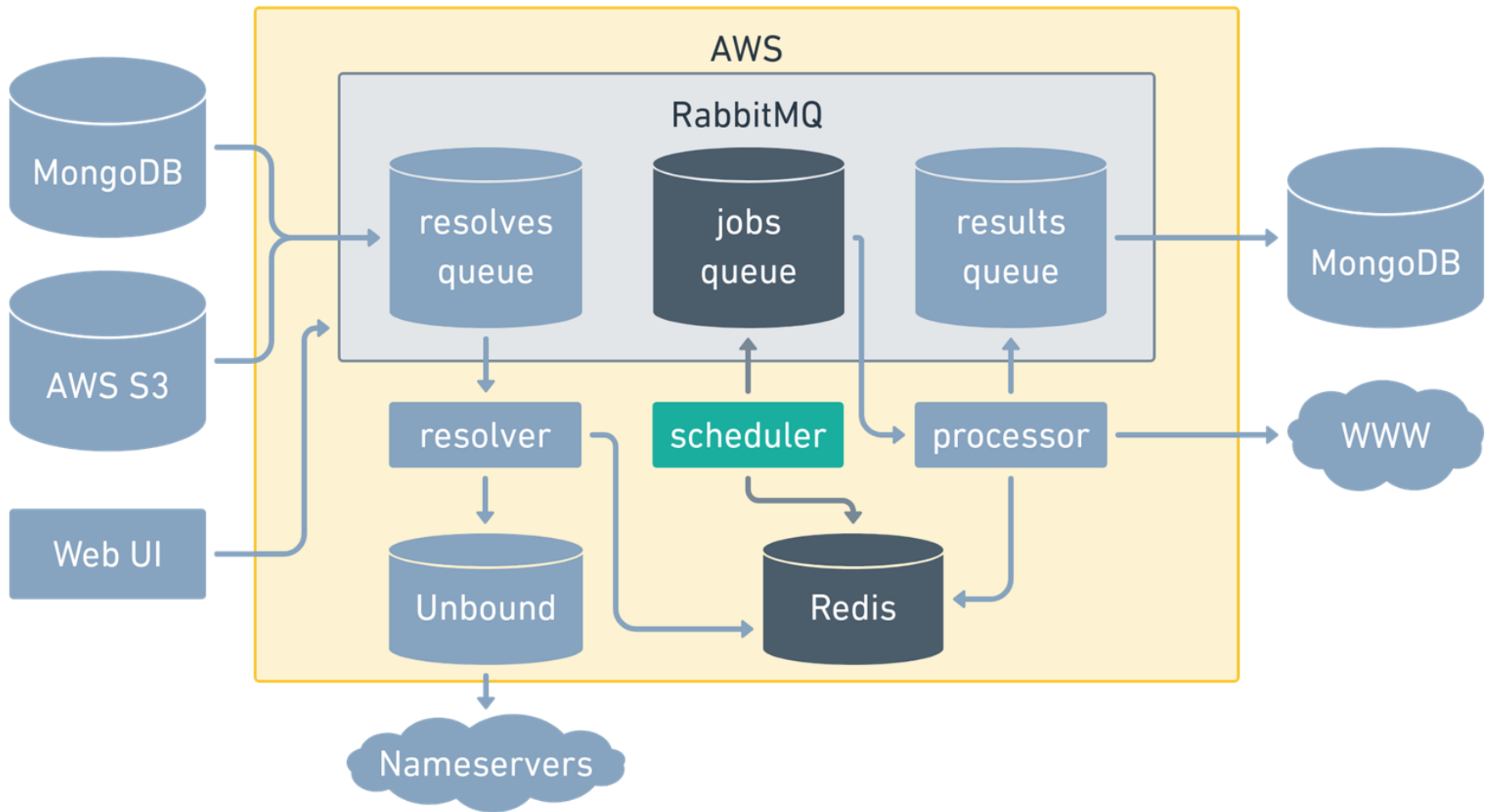


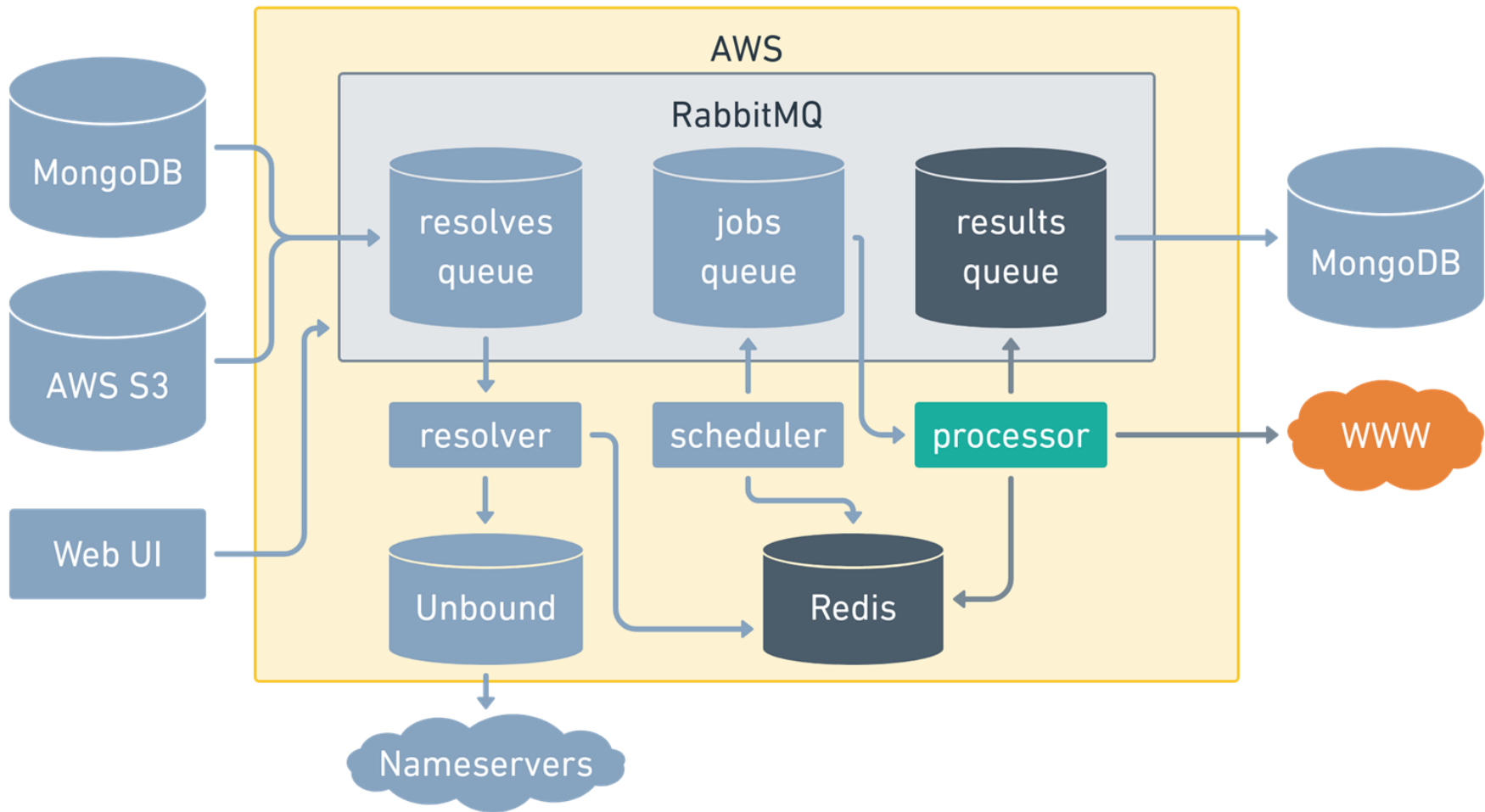
unbound



Много доменов на одном сервере

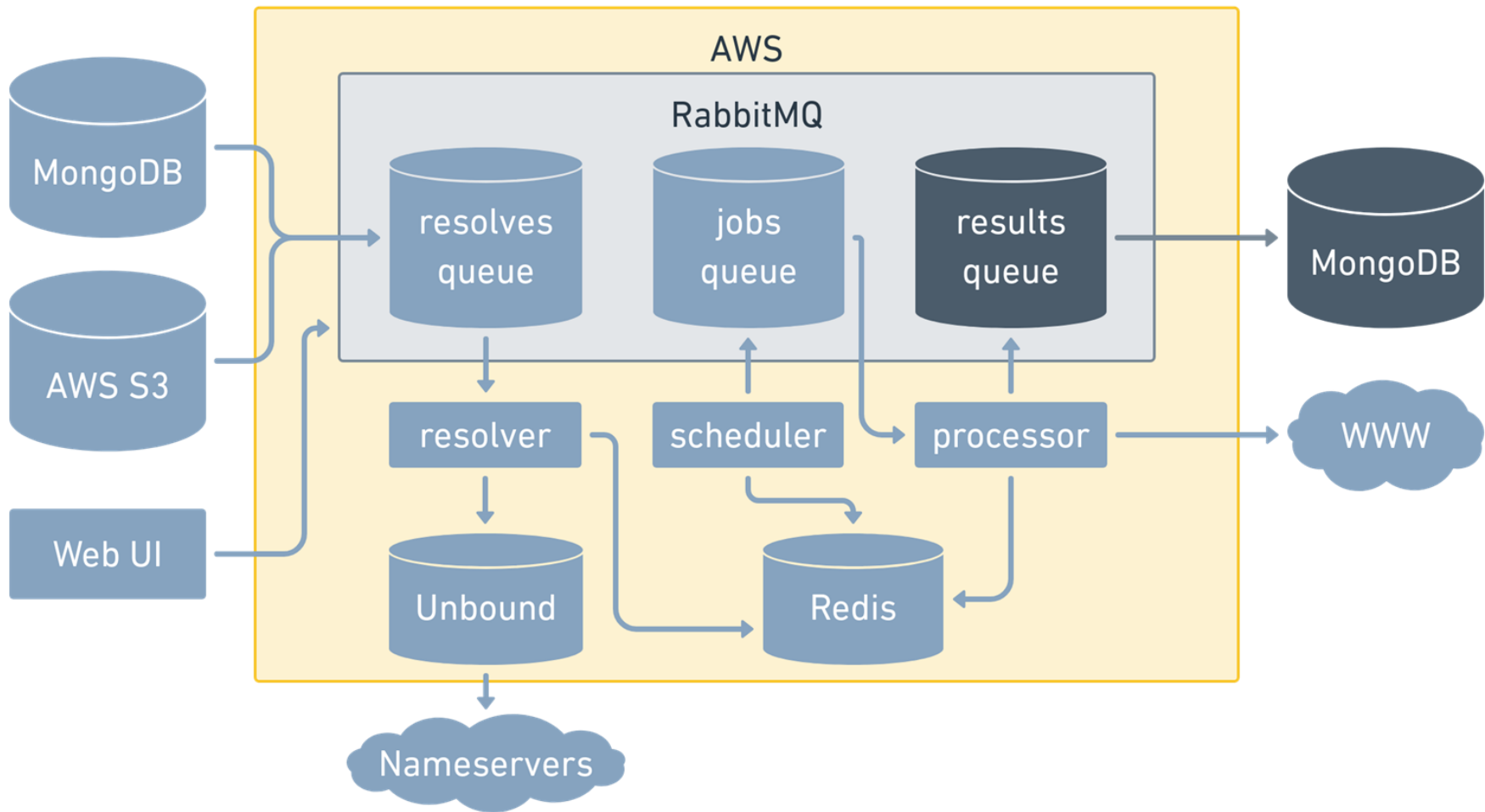
- Перемешивать список
- Делать лок по IP
- Делать паузы между запросами на один IP
- Используем ключи с ограниченным временем жизни в redis





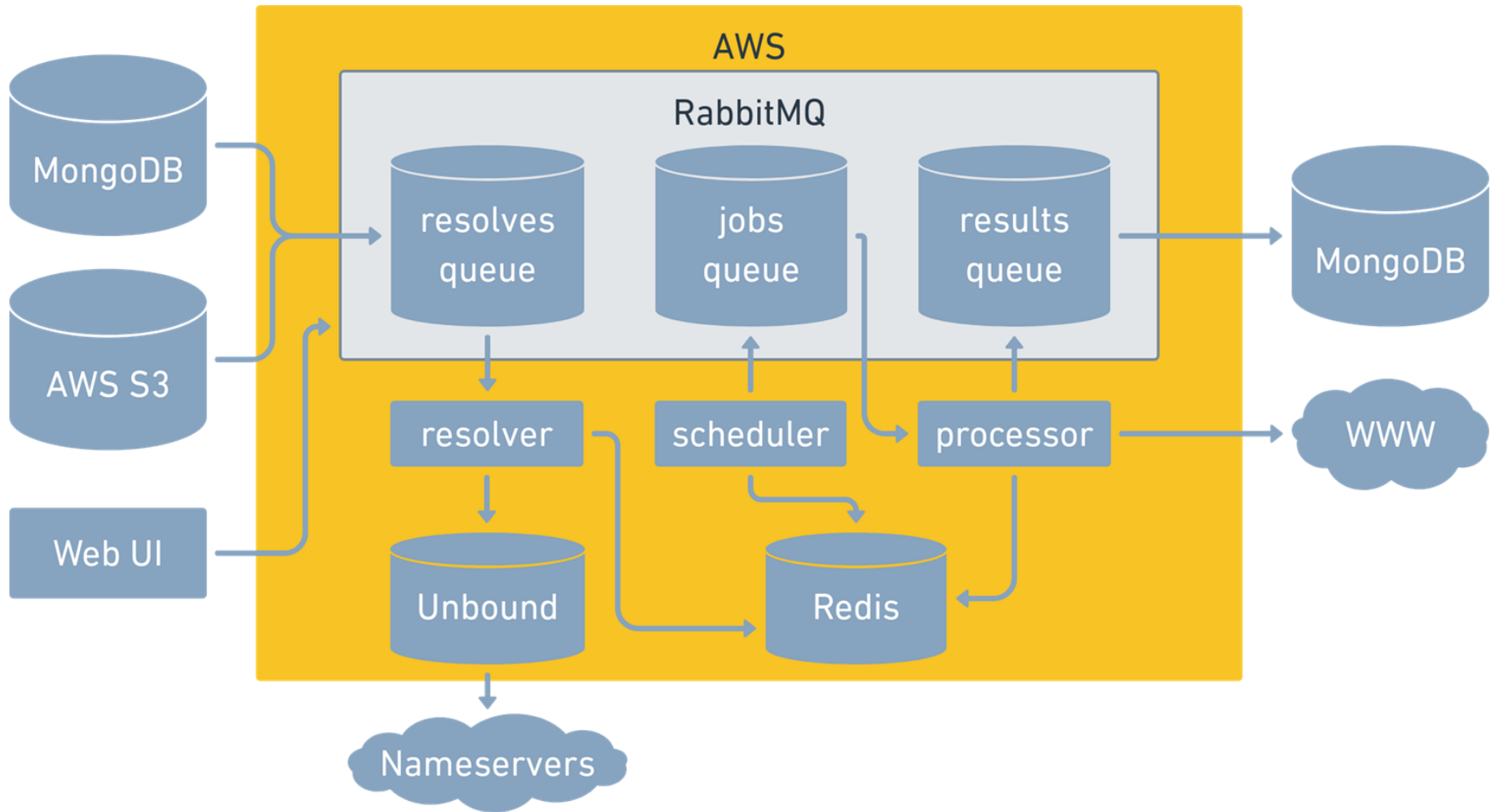
Масштабирование

- Используем ECS Fargate
- Автомасштабирование по CPU
- ≈ 100 -200 доменов в секунду на одном воркере



Сетевые задержки

- Перенесли всё (воркеры, RabbitMQ, Redis, Unbound) в один ECS-кластер
- Попробовали DocumentDB — спалили около двух тысяч долларов



Можно ли парсить сайты из AWS?

- Главное не наносить урон сайтам
- Поддержка амазона ориентирована на решение проблемы
- 3 инцидента за 4 года
- $\approx 300\$$ в месяц

Обработка большого объёма данных

- Streaming
- Batching
- Backpressure

Ограничение по памяти

- Node streams
- Async generators
- Backpressure
- Batching



Node streams

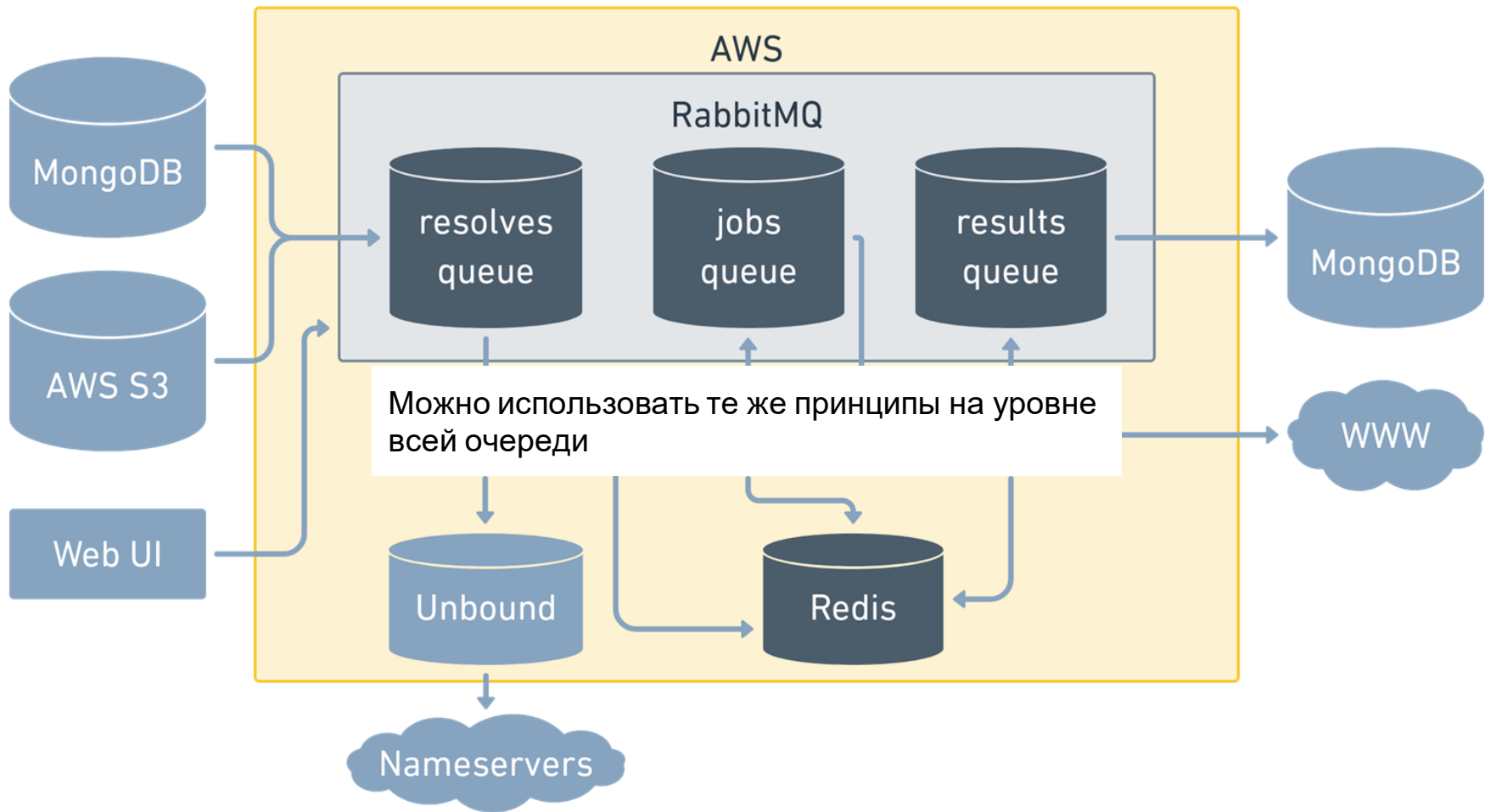
```
const stream = s3.getObject(params)
  .createReadStream()
  .pipe(csv.parse())
  .pipe(getFirstColumn)
  .pipe(createJob())
```

Async generators

```
async function* cursorToStream(cursor) {  
    while (await cursor.hasNext()) {  
        const doc = await cursor.next();  
        yield doc.domainName;  
    }  
}
```

Batching

```
const pool = new PromisePool(function* () {  
    for (const job of jobs) {  
        yield add(job);  
    }  
}, 500);  
await pool.start();
```



Дикий Дикий Веб

- Нельзя доверять входным данным
- Гигантские тайм-ауты в любой момент
- Гигантские ответы
- Плохие SSL-сертификаты

Устанавливаем лимиты

```
socket.setTimeout(timeout);
```

```
response.on('data', function (chunk) {
```

```
    body += chunk;
```

```
    if (body.length > bodyLimit) {
```

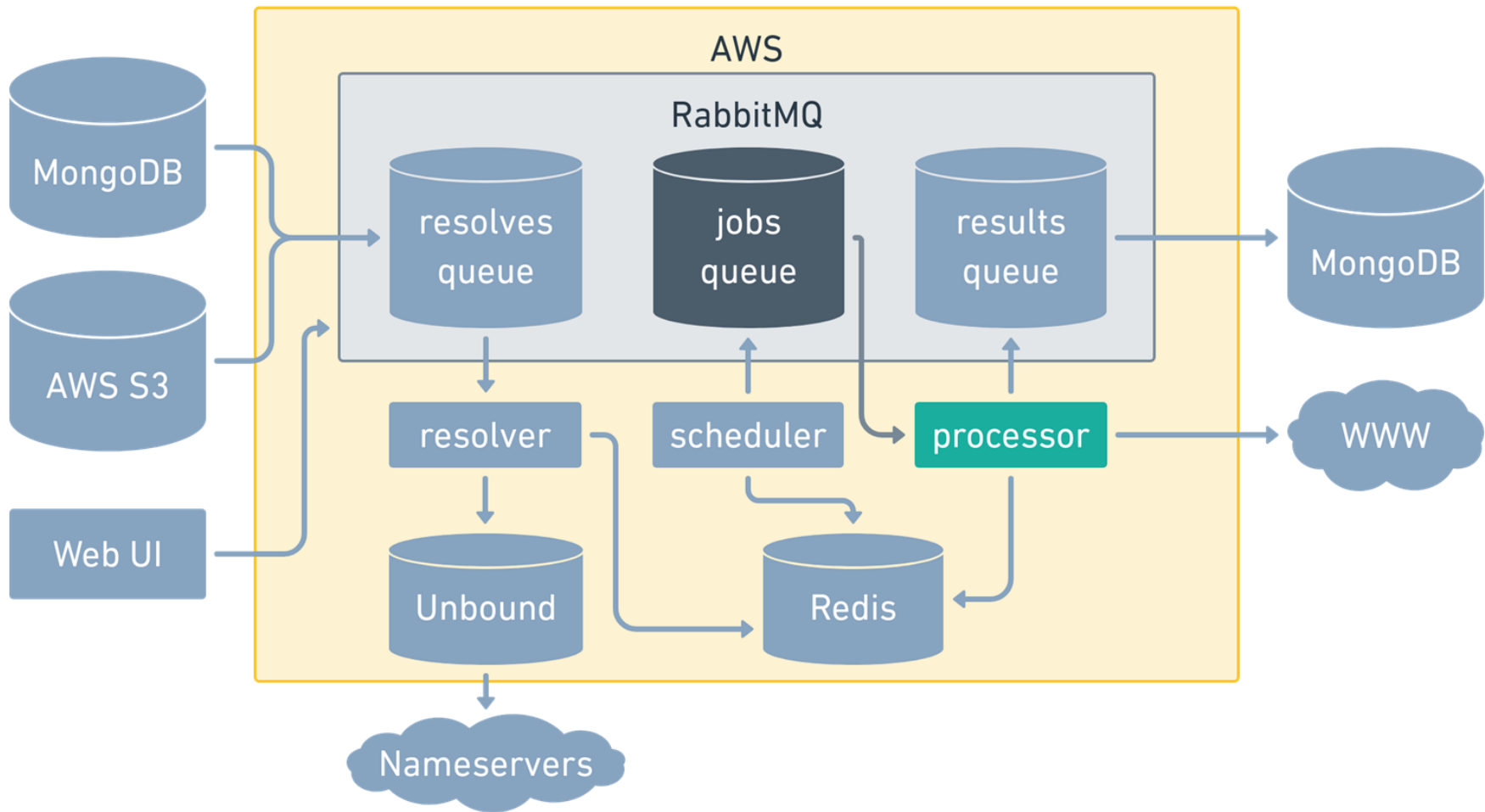
```
        request.abort();
```

```
    }
```

```
});
```

Плохие SSL-сертификаты

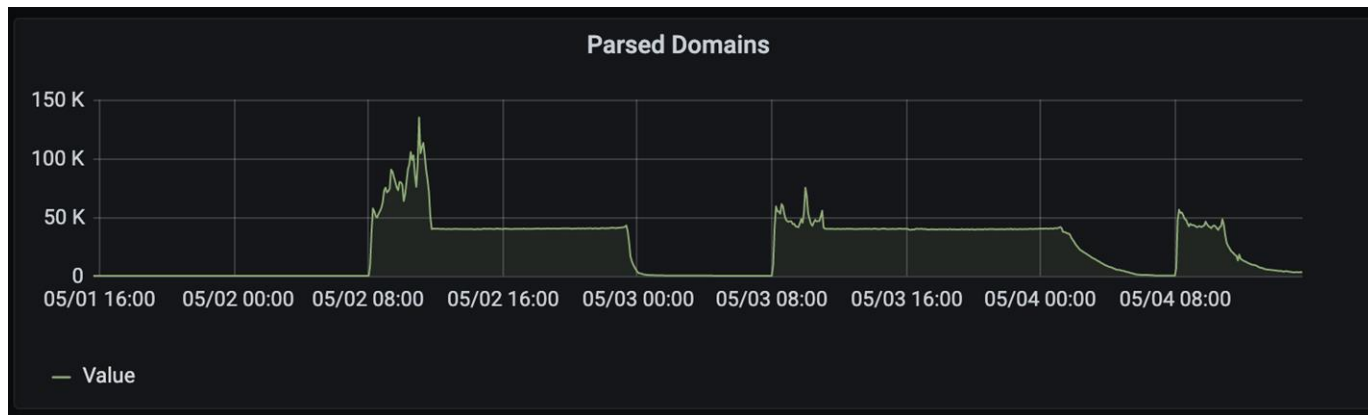
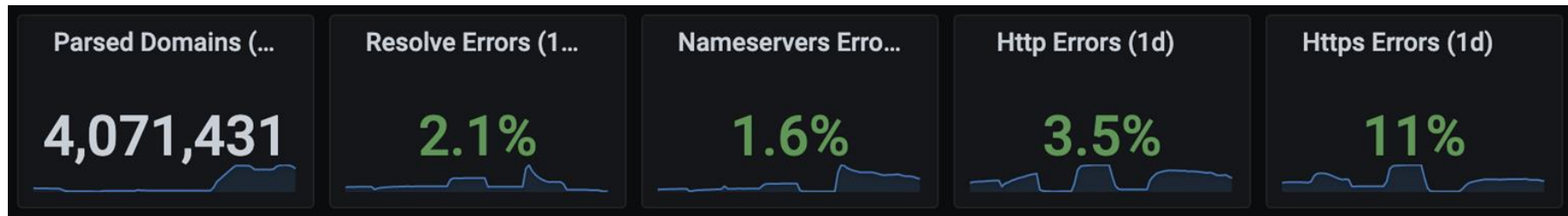
```
E0527 09:49:28.689000000 4712 ssl_transport_security.c:945]
  Handshake failed with fatal error SSL_ERROR_SSL:
  error:1000006b:SSL routines:OPENSSL_internal:BAD_ECC_CERT.
E0527 09:49:28.689000000 4712 handshake.c:241] Handshake failed
  with error TSI_PROTOCOL_FAILURE
E0527 09:49:28.689000000 4712 secure_channel_create.c:99] Secure
  handshake failed with error 1.
{ [Error] code: 14, metadata: Metadata { _internal_repr: {} } }
```



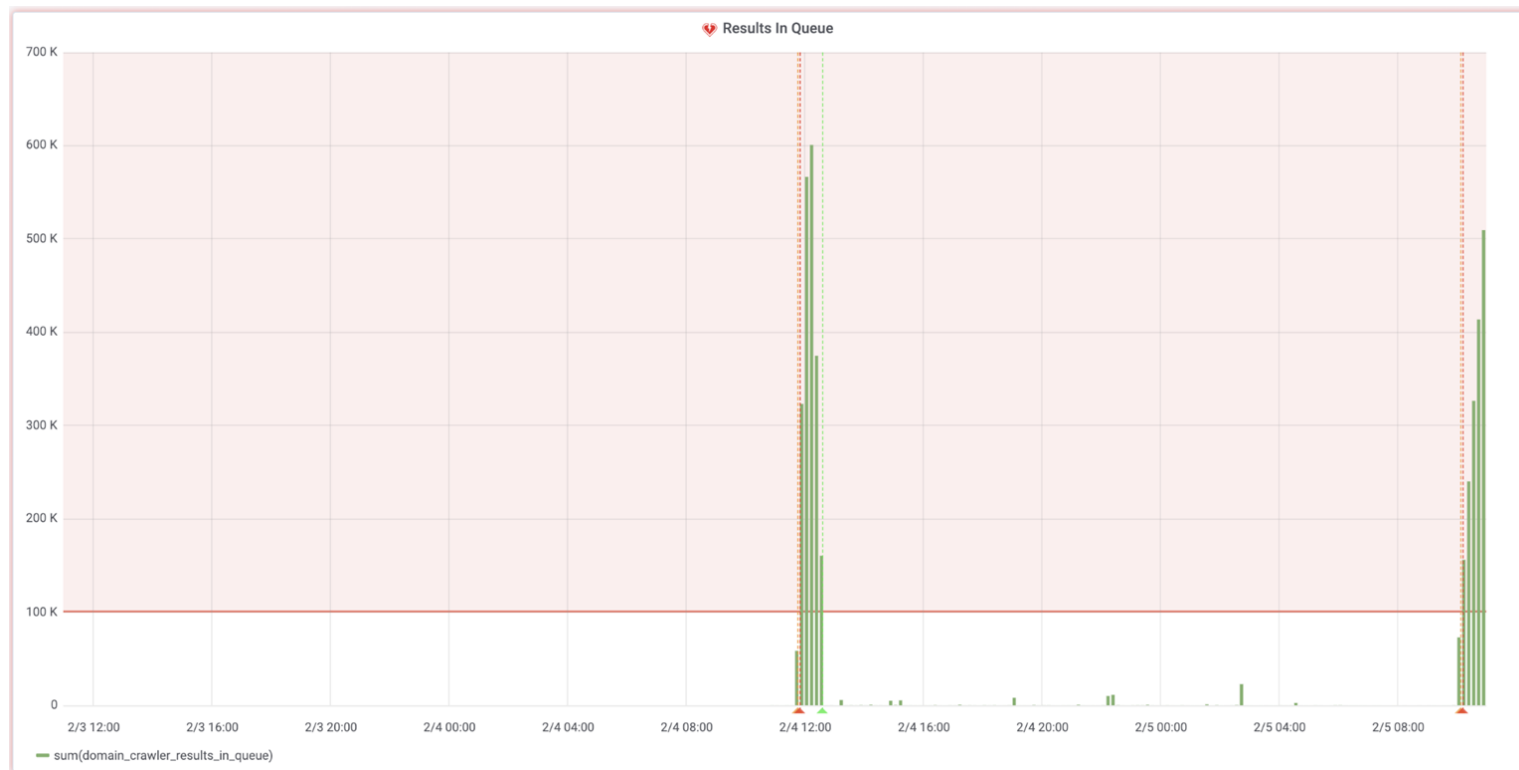
Если что-то может
пойти не так, оно
пойдёт не так...

...особенно если
нагрузка меняется на
несколько порядков.

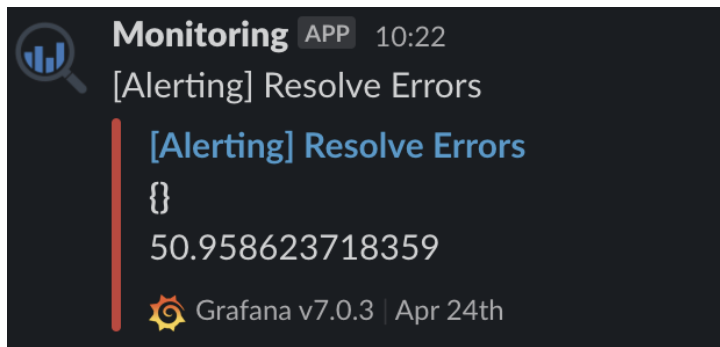
Мониторинг и оповещения



Мониторинг и оповещения



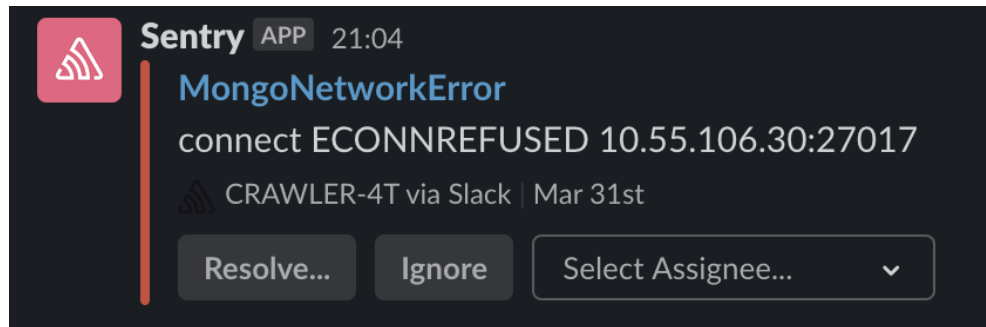
Мониторинг и оповещения



Monitoring APP 10:22
[Alerting] Resolve Errors
[Alerting] Resolve Errors
{ }
50.958623718359
Grafana v7.0.3 | Apr 24th



Monitoring APP 10:47
[Alerting] Jobs In Buffer
[Alerting] Jobs In Buffer
{ }
2105807.8571429
Grafana v7.0.3 | May 2nd



Sentry APP 21:04
MongoNetworkError
connect ECONNREFUSED 10.55.106.30:27017
CRAWLER-4T via Slack | Mar 31st
Resolve... Ignore Select Assignee... ▼

[ChatOps, избавляемся от рутины в RnD](#)

Итоги

- Боты должны быть “этичными”
- Можно использовать готовые наборы данных
- Нужно использовать сильные стороны инструментов
- Нельзя доверять ответам
- Автоматически отсеивать “плохие” домены
- Мониторинг и оповещения

Спасибо за внимание!

 ekaragodin@gmail.com

 <https://t.me/ekaragodin>

 <https://www.plesk.com/>

